

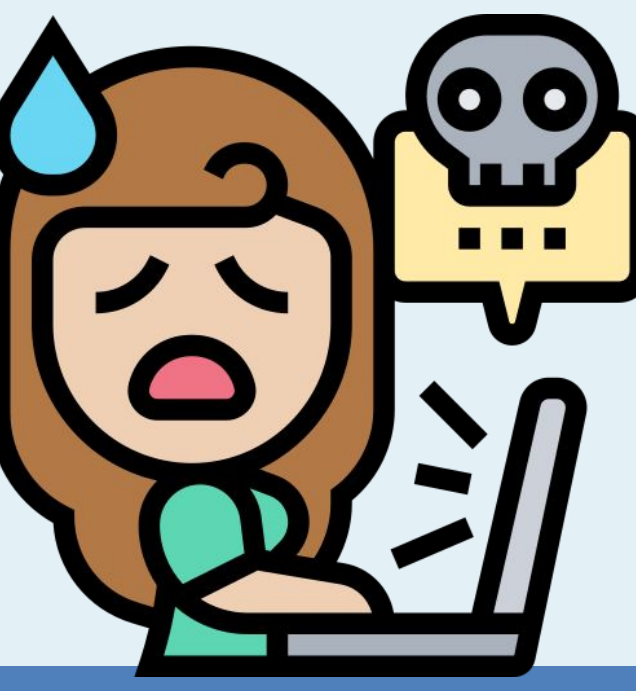


BiRDy: Bullying Role Detection in Multi-Party Chats

Anaïs Ollagnier¹, Elena Cabrio¹, Serena Villata¹, Sara Tonelli²

¹Université Côte d'Azur, Inria, CNRS, I3S

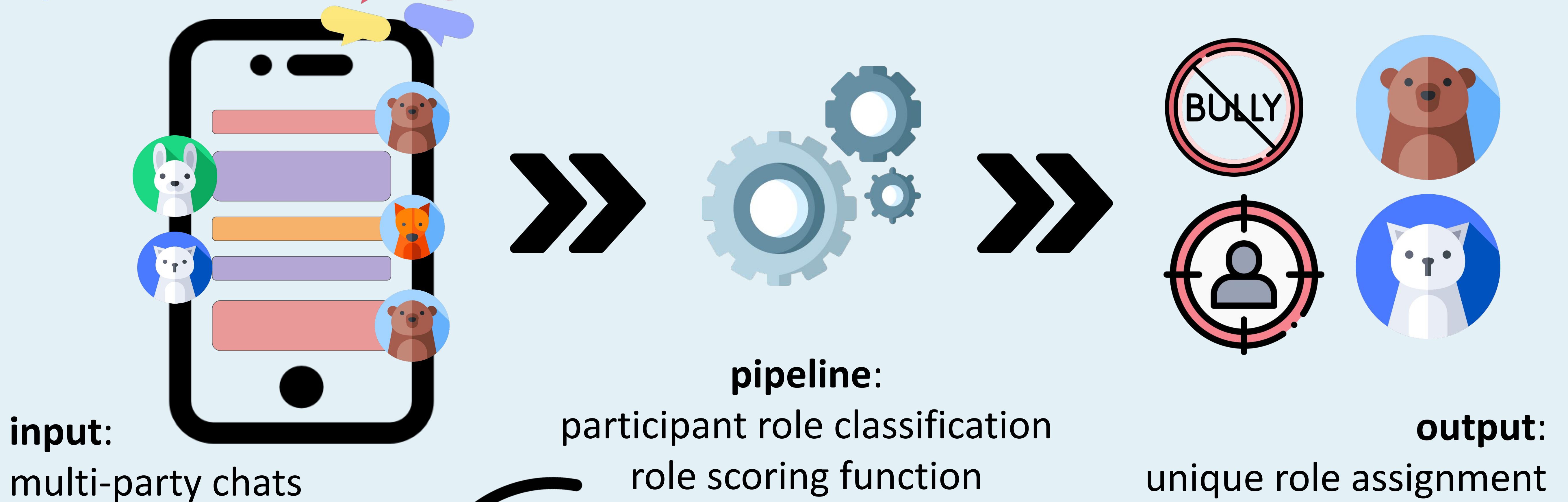
²Fondazione Bruno Kessler, Trento, Italy



BiRDy is a solution aiming at automating the identification of the different participant role of involvement in cyberbullying episodes, e.g. victim, bully, conciliator, bystanders assistant or defender

http://134.59.134.227/demo_prd/index.html

BirDy Overview



PRI classification

Preprocess:

Exploring the use of **contextual window** by collecting for each message the ***n* previous sentences authored by the same user**;

Modelling:

Classifying exchanged messages according to the role of the author (e.g. victim, bully, bystanders, conciliator) using **fine-tuned multilingual BERT models**.

Role scoring

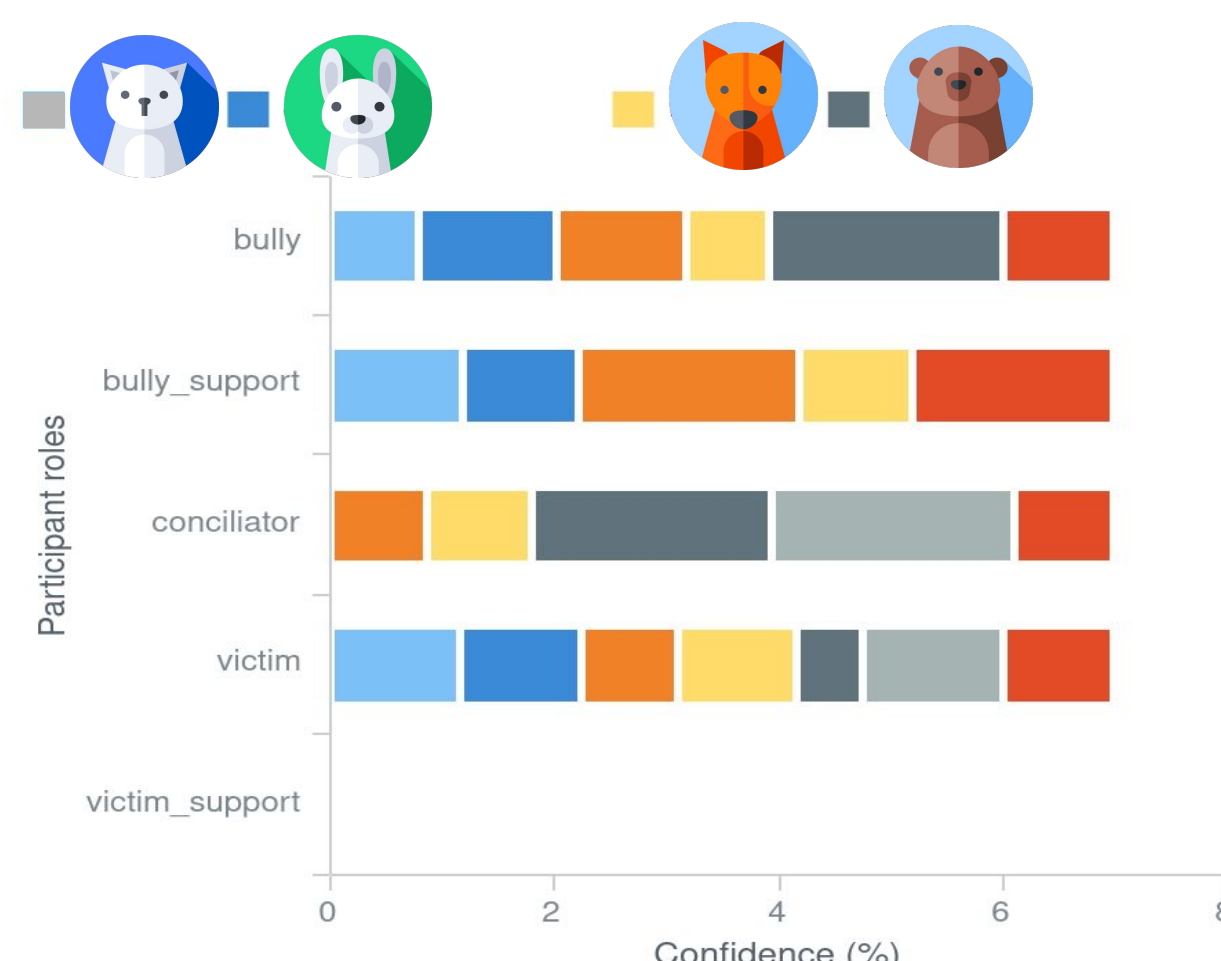
Assign a unique role to users by considering the probabilities established by the model.

$$Role_scoring(u, l) = \frac{\sum_{\substack{m \in speech_acts(u) \\ argmax(D(m))=l}} P_{am}(m, l)}{|speech_acts(u)|}$$

Confidence Score

Detailing for each label the **degree of likelihood** per user.

$$Confidence(u, l) = \frac{|U| \times Role_scoring(u, l)}{\sum_{u' \in U} Role_scoring(u', l)}$$

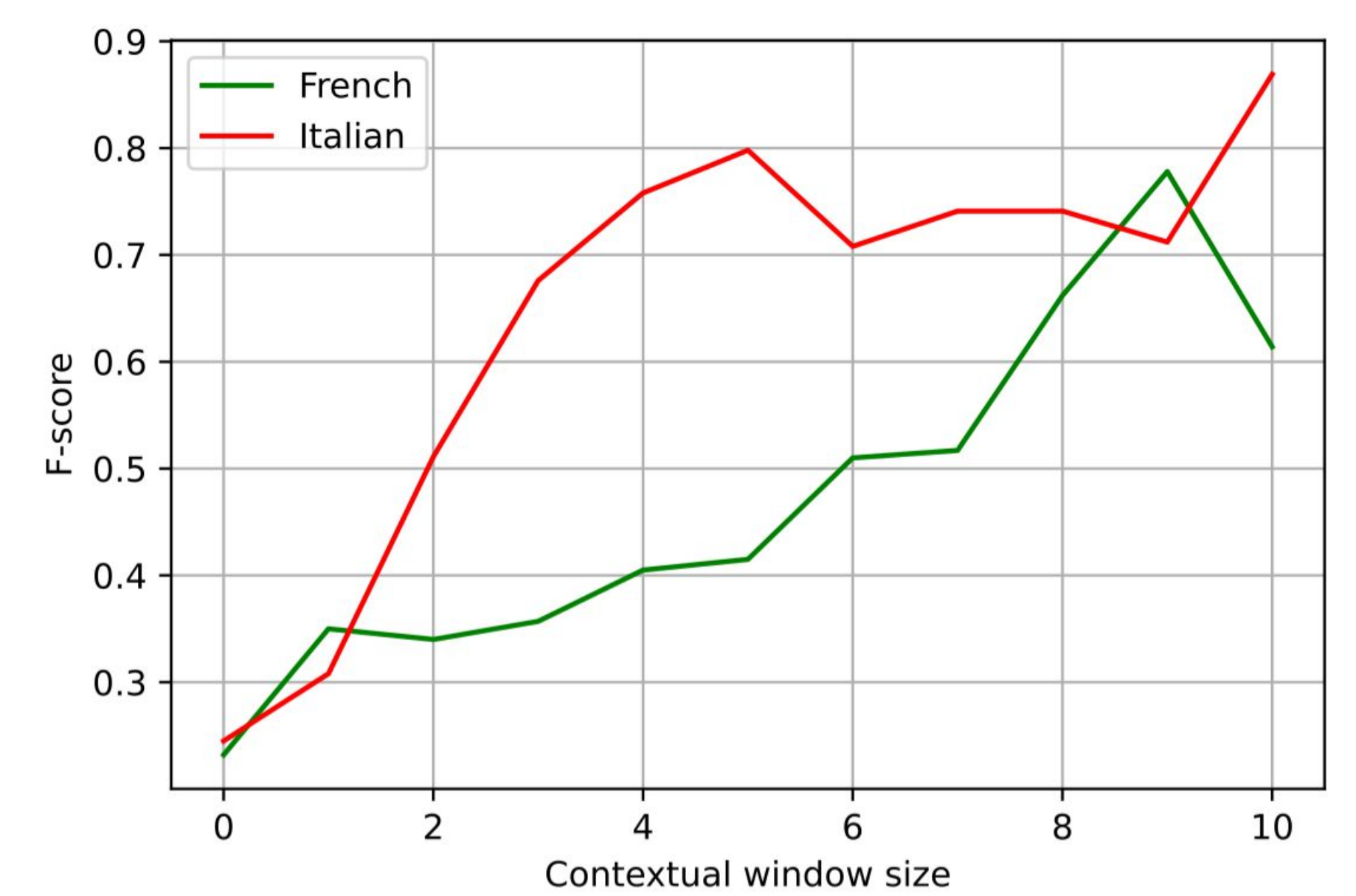


Evaluation

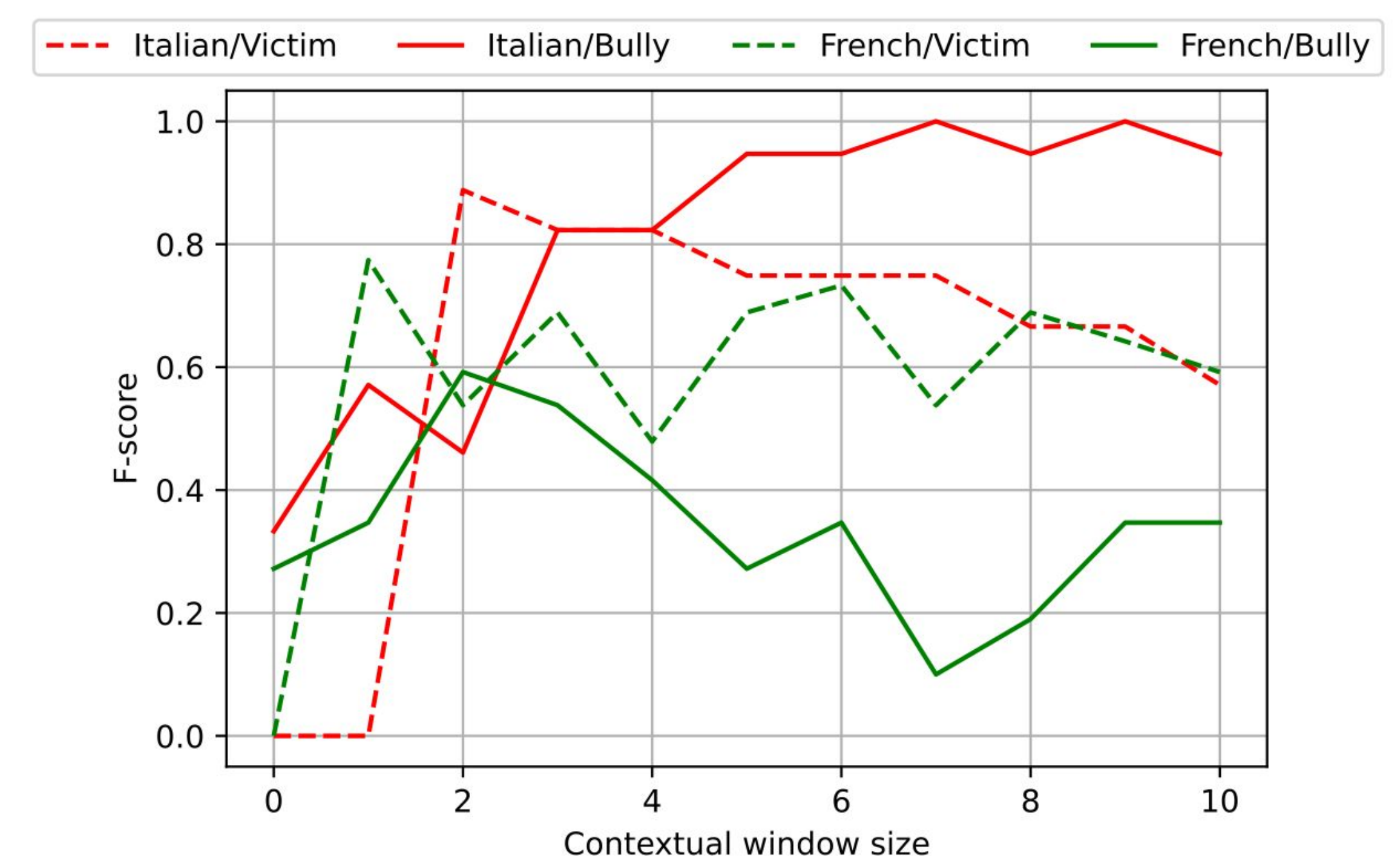
Datasets:

French: aollagnier/CyberAgressionAdo-v1
Italian: dhfbk/WhatsApp-Dataset

PRI classification



Bully/Victim classification



Acknowledgements

This work is supported by the French government, through the 3iA Côte d'Azur Investments in the Future project managed by the National Research Agency (ANR) with the reference number ANR-19-P3IA-0002 and the EFELIA Côte d'Azur project ANR-22-CMAS-0004.